

## 3rd International Conference on Industry 4.0 and Smart Manufacturing

## Smart extensions to regular cameras in the industrial environment

S. Brezani<sup>1,\*</sup>, R. Hrasko<sup>1</sup>, P. Vojtas<sup>1,2</sup><sup>1</sup>Globesys ltd, Framborska 58, Zilina, Slovak Republic<sup>2</sup>Dpt. Software Engineering, Charles University, Malostranske nam. 25, Prague, Czech Republic

---

**Abstract**

Data mining from unstructured data can be skillfully employed to improve the performance of manufacturing or industrial processes. The main goal of this paper is to create a fast emergency aid system for object detection in SME industrial premises. The basic assumption is that SMEs do not have any IT-trained personnel, and the solution has to be unsupervised edge computing.

We use several off-the-shelf models of deep neural networks pre-trained for smart city applications, ready for online object recognition and edge computing. Our system works without any retraining, additional annotation, or human intervention. Specifically, we present heuristics for the automated creation of PGT (Pseudo-Ground Truth). Based on PGT, we can automatically decide which model is the best in the specific environment. We present an application of fully automated enhancing image capture camera outputs to smarter ones. We evaluate our system in a controlled experiment. Low-resolution cameras and large areas cause problems for our method. We present a proof-of-concept for improving our system even in these challenging situations. The benefit is a knowledge extraction in a simple and inexpensive way to expand the organizations' databases with information from unstructured data from CCTV/IP cameras.

© 2022 The Authors. Published by Elsevier B.V.

This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0>)

Peer-review under responsibility of the scientific committee of the 3rd International Conference on Industry 4.0 and Smart Manufacturing

**Keywords:** Unsupervised learning; performance improvement of industrial processes; object recognition; pseudo-ground truth.; AI; machine learning in Industry

---

**1. Introduction with contributions**

Our primary interest in this paper is to use data mining on unstructured data to improve the performance of manufacturing or industrial processes. The main goal of this paper is to transform industrial premises images into

---

\* Corresponding author. Tel.: +421 911 802 025.

E-mail address: [samuel.brezani@globesys.sk](mailto:samuel.brezani@globesys.sk)

helpful knowledge for decision-making without retraining them for this specific environment. The main goal of this paper is to create a fast emergency aid system for object detection in SME industrial premises. The basic assumption is that SMEs do not have any IT-trained personnel, and the solution has to be unsupervised edge computing.

We use a variety of deep neural networks pre-trained in a smart city environment. Specifically, we use object detection to recognize several categories from regular camera images. The use-case is a customer using low-quality cameras on its premises. Moreover, we assume the customer needs fast help in extending the functionality of his camera system without any annotation, training, or human intervention.

Furthermore, we consider that customers have no skilled human potential for manual annotation of data. We present an application of fully automated enhancing image capture camera outputs to smarter ones. The benefit is a knowledge extraction simply and inexpensively to expand the organizations' databases with information from unstructured data from CCTV/IP cameras. Another significant trend is to optimize the energy efficiency and speed of object detection models. This trend makes it possible to use these models directly on the data source - edge computing. With this concept, it is possible to augment each CCTV/IP camera with smart features and thus enable data mining of unstructured data for Industry 4.0 needs.

The main contributions of our paper are:

- A new use-case in camera surveillance video object detection, asking for fast intervention.
- Selection of a variable system of models (pre-trained deep neural networks for object recognition tasks), mature enough to run online and as an IoT device.
- Fully automatic system enhancing a standard camera system with object recognition based on:
  - Heuristic system for unsupervised creation of pseudo-ground truth.
  - Automatically choosing the best model for each environment.
- Experiments evaluating the precision of our system:
  - Selection of representative data sources (video recording from different industrial environments).
  - Proof of concept of repairing weak results on "difficult" data, semantic knowledge, etc.

This paper is organized into 6 chapters. The first chapter describes the main contribution of this paper. The next 2 chapters introduced our novel method – define its prerequisites and describe the fully automated process creation of pseudo ground truth and model selection. The next chapter presents an evaluation of our method on human-annotated data. The final 2 chapters are dedicated to related work, conclusions, and our future work.

## 2. Videos with industrial premises and pre-trained models

Our goal is to use existing trained models of machine learning without any additional training to obtain "smart" information in the organization from their CCTV / IP cameras. We focus on ordinary cameras that do not support smart functionalities, respectively non "smart" software is supplied. The benefit of our solution for the organization could be a cheap and straightforward way to expand the organization's database with information from unstructured data from CCTV / IP cameras. The organization can use this data to control premises, parking lots, identify the use of premises, and trigger alarms, and the like. With more sophisticated analysis, an organization can identify trends in various events or use predictive analysis to predict future events. Some knowledge extraction should also be possible in the future.

### 2.1. YouTube videos with industrial environment

We focus on real examples of regular CCTV / IP cameras from industrial or business practice (for a few, see Fig. 1, more are depicted through the paper). For our research, we selected data from the youtube.com portal. We selected 13 videos from different environments (office, production line, parking, etc.) to obtain high variability.

Description of abbreviated names, full links, and video properties is listed later; see Table 3.



Fig. 1. Examples of regular images from CCTV / IP cameras (video source: 11\_Workspace, 12\_Alibi, 5\_Alibi).

## 2.2. Choice of deep neural networks *pre*-trained for object detection

A significant number of pre-trained computer learning models from the computer vision category are freely available on the Internet. Our primary interest is the area of industry. However, we did not find any training set from the industrial domain, e.g., office, production line, warehouse, parking lot, etc. Obtaining annotated data needed for training is a time-consuming activity that requires user intervention, and therefore, we want to avoid it. We chose several pre-trained models based on the COCO reference dataset [10]. This dataset contains more than 330,000 images containing annotations of 1.5 million instances of objects from 80 classes. But we can say that it is more from the area of "smart cities" and none of our target areas. Of course, we are aware that the accuracy of such models is probably significantly worse than in the case of their training in a specific environment (warehouse, office, etc.).

Many of these models use different architectures of deep neural networks. In our research, we selected models implementing the object detection task. This task identifies a set of instances of the objects in the image. For each object instance, the model returns:

- class - person, car, and the like,
- bounding box - position and size of the object - boundary where the object is located,
- confidence score with which the model predicts an object in a given bounding box.

To achieve better representability, we chose 9 deep neural networks based on different architectures, backbone layers and processed different sizes of input images for our research (see Table 1).

Table 1. Chosen models, backbone, and the number of identified object instances and classes for the confidence level.

Backbone	Confidence	cf > 0.3		cf > 0.5		cf > 0.7		cf > 0.9	
Backbone	Model / predicted	inst	cat	inst	cat	inst	cat	inst	cat
CSPDarknet53	YOLO4	13760	39	10366	27	7942	16	4486	10
ResNet50	Retinanet-ResNet50	23910	67	9800	36	4646	18	309	6
HourGlass104	CenterNet-HourGlass104	21939	60	9714	40	3476	16	185	3
ResNet50	CenterNet-ResNet50	14419	53	4527	20	849	8	9	1
EfficientNet B3	EfficientDet D3	20954	52	9883	27	4321	14	252	2
EfficientNet B6	EfficientDet D6	27556	47	11507	31	5807	14	344	3
ResNet152	Faster R-CNN	19947	70	12026	58	8343	46	4299	18
ResNet101	Retinanet-ResNet101	23712	62	6953	35	1275	7	7	3
Darknet53	YOLO3	13935	42	9344	34	7118	21	4802	11

### 3. Fully automated creation of pseudo-ground truth

#### 3.1. Model analysis

For our needs, we extracted images from the video (frames with a frequency of 1 per sec). This parameter is discussed later – it can influence the quality of our service. Using all models, we performed predictions on each of the 3838 images (see Table 3 on video length). Table 1 shows model ID with backbone architecture and the number of identified instances and classes of objects for each model at different confidence levels (sum on all videos). Note that we are not talking about correctness - we take the position of correctly marking classes and instances based on the observed legalities – to have a fully automated solution.

Nevertheless, we can see that the simple use of neural networks and simple voting is not working. Networks trained on smart city data recognize too many classes and instances that clearly do not fit into the industrial environment (as proof of concept, we have experimented with semantic methods to recognize classes that do not fit into the premise).

#### 3.2. Our solution – automated creation of pseudo ground truth

An important criterion for our method is to minimize the need for human intervention and support the automation of such a process. Based on the established criteria, we proposed a method for identifying the most suitable machine learning model from N available models. The frame is classified by all models, and we design a "merge prediction" method deciding which instances with classes will be chosen to our ground truth. We extract images from the input video at the selected frequency (in our case, 1 second). For each image, we make a prediction using each model. This way, we get a set of identified objects from all models from the input image. For further processing, we select only those predictions whose confidence > 0.5.

Because different models assign different bounding boxes to an instance, our heuristic assumption is (if correct) that these bounding boxes are not very different. The next step is to combine these predictions using our algorithm to obtain a "pseudo ground truth". We call it "pseudo" because it is created fully algorithmically without human intervention during annotation. The algorithm for combining predictions to form "pseudo ground truth" is something like this:

```

pairs_queue = []
bbox_group = {}
FOR bbox_a IN bboxes:
    FOR bbox_b IN bboxes:
        IF bbox_a != bbox_b && iou(bbox_a, bbox_b) > threshold: pairs_queue += (bbox_a, bbox_b, IOU(bbox_a, bbox_b))
SORT(pairs_queue, KEY=2, DESCENDING) # descending sort by iou values
FOR idx, bbox IN ENUMERATE(bboxes): bbox_group[bbox_a] = idx
FOR [bbox_a, bbox_b, iou_score] IN pairs_queue: ga, gb = bbox_groups[bbox_a], bbox_groups[bbox_b]
    IF ga == gb: CONTINUE
    merge_candidates = SET([bbox IN bboxes IF bbox_groups[bbox] IN (ga, gb)])
    model_candidates = SET([MODEL(bbox) FOR bbox IN merge_candidates])
    IF COUNT(merge_candidates) != COUNT(model_candidates): CONTINUE
FOR bbox in merge_candidates: bbox_group[bbox] = ga

```

Algorithm inputs are all frames from a video. For each frame separately, the algorithm searches for all prediction pairs whose IOU - intersection over union metric is greater than the threshold parameter (in our case, 0.5). It stores such pairs in the pairs queue field, sorts these pairs in descending order according to their IOU value. Then the algorithm initializes a key variable `bbox_group` representing each prediction assignment. Subsequently, the algorithm iterates through the pair of predictions. Let us comment here that this algorithm is a pure heuristic. We do not aim to have a correct algorithm because we do not have annotated ground truth. Instead, we aim to create a pseudo-ground truth formed by model candidates automatically.

In the next step, we initialize the set of all predictions that belong to the group. These are merge candidates – elements of future pseudo ground truth. Each group has 1 to N predictions (N is the number of models), each coming from a different model. Because only those predictions whose  $IOU > \text{threshold}$  (which we consider to be a sufficient condition for "identity" of objects) are combined into groups, we consider all predictions in a group to be the identification of the same object.

Fig. 2 shows an example of combining predictions that we achieved with the algorithm. The image on the left shows color-coded predictions from different models represent how many of the models identified each object. The right part of the figure shows the already associated predictions together with the score (number of models agreeing). Fig. 3 shows a more complicated situation in which two people overlap. Our algorithm, even in this case, combined the predictions correctly.

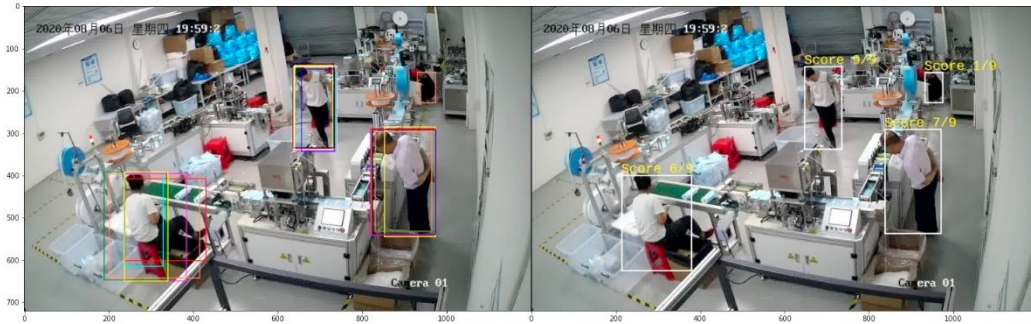


Fig. 2. Left bounding boxes (bbx) of various models, right bbx with a number of models which object recognition agrees – these form an algorithmically created "pseudo ground truth" without any human intervention (video source: 4\_APL80).

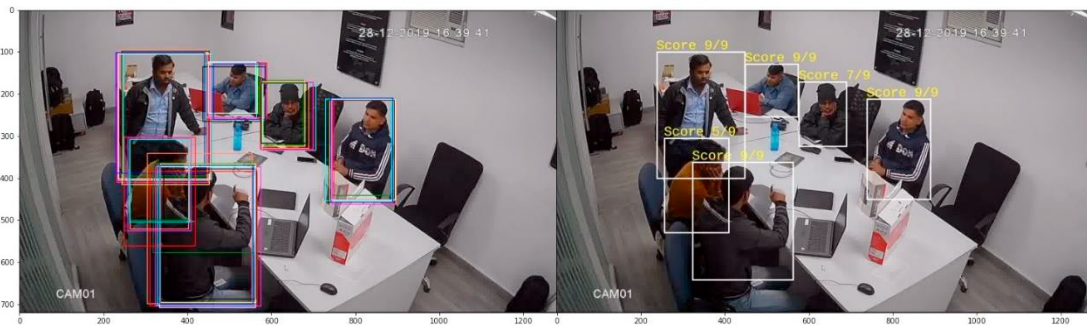


Fig. 3. More complicated situation in which two people overlap. Even in this case, the algorithm combined the predictions correctly (video source: 0\_CCTV).

After processing the image according to the algorithm, we obtain grouped predictions for each image. We call them pseudo-ground truth. For each such prediction, we know how many of all the models identified a given object. At this point, we define a key assumption: if there is a real object in the image, then we assume that at least  $K$  of the  $N$  models identifies this object. In our research, we worked with  $N = 9$  models and chose  $K = \{3, 5, 7\}$ . It should be noted that in this heuristic, we do not consider at all the confidence with which the individual models recognize the instances - and the groups could also be assigned "group" confidence. Using this assumption, we obtain a set of pseudo ground truth ( $PGT_K$ ) objects for each image and  $K = \{3, 5, 7\}$ .

### 3.3. Choosing the best model by the highest recall

We already have all the necessary input data - images and identified objects to select the most suitable model. For each  $K = \{3, 5, 7\}$  and each model, we identify the objects in the images and compare them to  $PGT_K$ . The main point of our selection of the best model  $m_K$  is that we chose the one with the highest recall (over all frames, measure see [12]). We find this an important part of our heuristics because we would like our customers to have rather higher recall than higher precision. This also lowers the false negatives rate.

To our surprise, the best results were obtained for  $K = 3$  (in the sense that recall was higher than for  $K = 5, 7$ ). Nevertheless, we wanted to evaluate our method more precisely; see next chapter.

#### 4. Validation

To validate the accuracy/precision of our heuristics, we manually annotated a portion of the input data for each video. We then performed object identifications for the annotated data by all models with a confidence of at least 0.5. Subsequently, we compiled a list of models for each video, sorted by average precision ([12]) over all frames. Supervised object detection tasks can be evaluated by multiple metrics (see [13]). We chose average precision for its simple human interpretation. For comparison, we also evaluated these metrics for the models recommended by our heuristic and human-annotated GT.

Table 2. Comparison of automated PGT versus human-annotated GT and best model metrics.

Video source	Computed PGT <sub>3</sub> / human-annotated GT		The best model m <sub>3</sub> performance on human-annotated GT		
	Average Precision	recall	Average Precision	recall	Best chosen by our system
0_CCTV	0.9944	0.9944	1.0000	1.0000	YOLO3
10_MAGGI	0.8235	0.8235	0.9108	0.9118	EfficientDet D6
<b>11_Workspace</b>	<b>0.1856</b>	<b>0.1926</b>	<b>0.2664</b>	<b>0.2704</b>	<b>EfficientDet D6</b>
12_Alibi	0.6984	0.6984	0.6974	0.6984	EfficientDet D6
<b>13_SLEEPING</b>	<b>0.2000</b>	<b>0.2000</b>	<b>0.2000</b>	<b>0.2000</b>	<b>YOLO4</b>
<b>2_4K</b>	<b>0.3329</b>	<b>0.3425</b>	<b>0.3650</b>	<b>0.3836</b>	<b>EfficientDet D6</b>
3_BMW	1.0000	1.0000	1.0000	1.0000	Faster R-CNN
4_APL80	0.8199	0.8288	0.6089	0.7055	EfficientDet D6
5_Alibi	0.9592	0.9592	0.9592	0.9592	Faster R-CNN
6_HD	1.0000	1.0000	0.9133	0.9286	Retinanet-ResNet50
7_Allied	0.6452	0.6452	0.5914	0.5914	EfficientDet D6
8_Warehouse	0.5000	0.5000	0.5000	0.5000	EfficientDet D6
9_MESSOA	0.8115	0.8333	0.7976	0.8333	CenterNet-HourGlass104
0_CCTV	0.9944	0.9944	1.0000	1.0000	YOLO3

The bold rows in Table 2 display data sources with a low average precision value (<0.4). Looking at the videos' quality attributes, we can get hints where our method is weak depending on video properties, see Table 3.

Table 3. List of video properties used for evaluation later. See Appendix for more video information.

Short name	Premise	Frames	Resolution	Area size	Number of persons	Link
0_CCTV	office	31	high	small	few	<a href="https://youtu.be/mMXDU8q4fIQ">https://youtu.be/mMXDU8q4fIQ</a>
10_MAGGI	manufacturing	174	high	small	few	<a href="https://youtu.be/cIVDZiJ7NjY">https://youtu.be/cIVDZiJ7NjY</a>
11_Workspace	manufacturing	182	low	medium	many	<a href="https://youtu.be/xBZzc8iZjT8">https://youtu.be/xBZzc8iZjT8</a>
12_Alibi	warehouse	55	high	medium	few	<a href="https://youtu.be/O2pG39jO0JE">https://youtu.be/O2pG39jO0JE</a>
13_SLEEPING	warehouse	34	low	small	few	<a href="https://youtu.be/ZdZeja401xM">https://youtu.be/ZdZeja401xM</a>
2_4K	parking	275	high	large	many	<a href="https://youtu.be/wqlO5i3N-FU">https://youtu.be/wqlO5i3N-FU</a>
3_BMW	manufacturing	858	high	small	few	<a href="https://youtu.be/P7fi4hP_y80">https://youtu.be/P7fi4hP_y80</a>
4_APL80	manufacturing	1409	low	medium	many	<a href="https://youtu.be/yZcufAb7sMk">https://youtu.be/yZcufAb7sMk</a>
5_Alibi	warehouse	29	high	small	few	<a href="https://youtu.be/0xHvoobqHhE">https://youtu.be/0xHvoobqHhE</a>
6_HD	office	193	high	small	few	<a href="https://youtu.be/O6IMyzDEMHE">https://youtu.be/O6IMyzDEMHE</a>
7_Allied	warehouse	178	high	large	few	<a href="https://youtu.be/CxpRrRBxl4Q">https://youtu.be/CxpRrRBxl4Q</a>
8_Warehouse	warehouse	265	high	small	few	<a href="https://youtu.be/YQtBjpReiaM">https://youtu.be/YQtBjpReiaM</a>
9_MESSOA	warehouse	155	high	small	few	<a href="https://youtu.be/OO7XT24AmTY">https://youtu.be/OO7XT24AmTY</a>
	Total	3838				

The area size and resolution are defined in Table 3; we can visualize these data sources using the graph in Fig. 4. We can see the graph's dependency between resolution, scene area size, and average precision value. A more detailed analysis is a goal in our further research. Some initial experiments are described in the Conclusions.

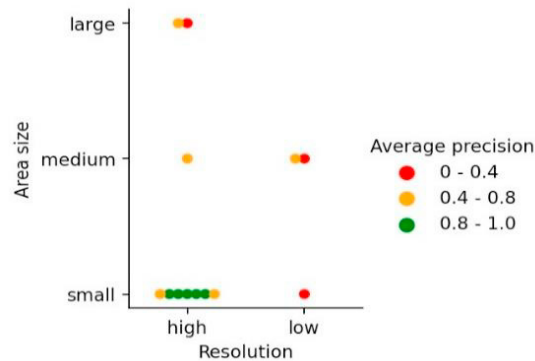


Fig. 4. The dependency between video attributes and the average precision value of the models selected by our heuristic.

## 5. Related work

We can divide related work into 3 parts. The first part is focused on object detection because our methods benefit from the variability of these methods. The second part describes a relation between our methods and existing methods of models selection. The last part deals with applying our method to the context of IoT, data mining, and knowledge extraction.

**Object detection** methods have been the subject of research for more than 20 years, see, e.g. [19], and are now long established in industrial practice (such as face detection, pedestrian detection, traffic sign detection, text detection, etc.). Research in this area has focused on several directions.

The first is improving object detection models themselves. One of the first architecture is R-CNN [7]. This architecture was one of the first deep neural network methods for the task of object detection. Other architectures such as [8] and [17] improve the original R-CNN architecture, succeeding in significantly increasing the speed of the model. Another breakthrough architecture was [16], which was the first to use the concept of one-step object detection. This concept brought significantly higher model speed but at the expense of accuracy. Other architectures such as [1] and [18] improved an original YOLO architecture while optimizing speed and energy efficiency. Optimizing object detection models' energy efficiency and speed allows them to be used directly at the data source (edge computing) and brings the ability to connect these devices to IoT networks. Furthermore, using this concept makes it possible to extend each CCTV/IP camera with smart features and thus enable data mining of unstructured data for Industry 4.0 needs. Nevertheless, they all need an annotated data set to be trained.

Another research direction focuses on specific improvements to selected problems of existing methods. For example, the method [2] focuses on minimizing false positive detections of nearby or overlapping objects. Other methods focus on eliminating hard (with very large confidence) false positives, see, e.g. [4] and [3]. These methods are based on extending the detection model with an additional classification model. We also identified problems with hard false positive and false-positive detections of nearby and overlapping objects during our research. However, in our case, we could not use these methods because they are from the supervised learning category.

**Model selection** and its evaluation are most often accomplished by training and evaluating the model on disjunctive datasets - the training and validation sets. These sets are generated using methods such as K-Fold, LOO, LPO, and others ([14]). In our case, we cannot use these methods directly because we need the ground truth to evaluate the model using them.

Further research in the field of AutoML has also brought the possibility of automatic model selection. Framework Auto-Keras [9] uses Bayesian optimization to select the most appropriate neural network architecture.

Framework Auto-Sklearn [5] uses various meta-learning techniques not only to identify the best model but also to build a pipeline for data preprocessing. Similar to the previous case, we cannot use these methods directly because



they are from the supervised learning category as they require ground truth data. As can be seen, for selecting the best model, a set of ground truth data is required. Our method introduces a heuristic by automatically generating a set of "pseudo ground truth" data.

Our method uses a similar idea to learning from crowds methods as [15], just our crowds is a crowd of deep neural networks. These methods are based on the principle that multiple annotators - experts with different experiences annotate the same data. In this way, we obtain multiple subjective sets of ground truth data—subsequently, these methods approximate ground truth by different approaches - majority voting, maximum likelihood estimation etc.

**Method application:** In [11], the authors show how to utilize data mining techniques which are now within reach of numerous smaller manufacturing operations. Moreover, they provide a further understanding of how moves towards fully Industry 4.0 ready factories may be made in the years to come.

In [6], the authors review existing data mining (both unsupervised and supervised) and analytics applications in the process industry over the past several decades. They highlighted and discussed several perspectives for future researches on data mining and analytics in the process industry.

A basic reference for MDM-multimedia data mining is [21]. Our research falls into Part II, "Multimedia Data Exploration and Visualization", especially chapter 5. "Mining Rare and Frequent Events in Multi-camera Surveillance Video unsupervised [26]", several other chapters are relevant, and, e.g. [27] and can serve for later comparison of results.

Nevertheless, none of these resemble our goal. Concerning novelty, we did not find any mention of such or similar use-case situations in the literature. Another step of Disruptive Innovation in the MDM area reflects shifts toward networked social communities, mobile devices, and sensors (see [23]). Some use of deep neural network techniques in MDM is described in [28]. However, [28] claim that one drawback is that these methods mostly used supervised approaches where a large corpus of labeled training data is needed. Whereas some studies used unsupervised DL methods, we expect to see more research in this area in the near future. We hope that our results contribute to this expectation, and the PGT idea will find wider use in experiments.

Authors of [20] propose EdgeMiningSim, a simulation-driven methodology inspired by software engineering principles for enabling IoT Data Mining. In our case, CCTV/IP can be considered as IoT devices. The PGT creation can be run on-premise or in a cloud (as AI as service), but the final chosen model runs at the customer's site. Therefore, appropriate modification of [20] can have an impact on our IoT data mining.

With an increased number of our system customers, we can use knowledge-based video indexing and content management framework for domain-specific videos from [22]. Moreover, for a new customer via associations (similarity) to previous cases, this can help us to choose the right DNNs models for new PGT creation.

We would like to emphasize that we use DNN models without any retraining, just as they are on the self. With longer use and possible feedback from the user, we can think of using fine-tuning and transfer learning as described, e.g., in [24] and [25].

In [29], the authors discuss the problem with some deep learning architectures that display problematic behaviors, such as confidently classifying unrecognizable images as belonging to a familiar category of ordinary images and misclassifying minuscule perturbations of correctly classified images. All the examples they mention are artificially generated either by white noise or evolutionary computing.

[30] is a recently edited volume containing an advanced report on current problems concerning multimedia in the IoT. Another good source of current development is also [31].

Let us mention at the end that none of these resembles our goal: to create a fast emergency aid system for object detection in industrial SMEs. Our basic assumption is that SMEs do not have any IT-trained personal and the solution has to be unsupervised edge computing working on "non-smart" cameras.

## 6. Conclusions, future work

The main goal of this paper is to create a fast emergency aid system for object detection in SME industrial premises. The basic assumption is that SMEs do not have any IT-trained personnel and the solution has to be unsupervised edge computing. We use several deep neural networks pre-trained for smart city applications, ready for online object recognition and edge computing. Our system works without any retraining, additional annotation, or human intervention. Our research has provided a method for selecting an object detection model for a particular data source.



In our method, we used a simple mechanism - based on the principle of learning from crowds to generate a set of pseudo ground truth data. We expect that by using more sophisticated methods for generating pseudo ground truth, the accuracy of selecting the best model can be improved. We hope the benefit is a knowledge extraction in a simple and inexpensive way to expand the organizations' databases with information from unstructured data from CCTV/IP cameras. Thanks to the new paradigm of Industry 4.0 and by employing information technology, this process becomes easier.

A challenge remains false negative. Optimizing our method on recall lessens this problem. In addition, an interactive deployment can help. E.g., giving a permanent operator the possibility to comment on our output (just clicking what is correct-true positive, what not-false positive, what is missing-false negative, etc.).

We have to have prepared default models for initial functioning from time  $t = 0$  on, for various industrial environments. Simultaneously, we start to record a video stream. Then, we select a time for learning  $t_1$  where a sufficient period with typical operation repeatedly occurred. Let us run the above algorithm on these records and select  $K$  with the best result on  $PGT_K$ . Then, we select the best model,  $m_K$ , and start to use it for this video.

We intend to extend our method by continuously evaluating model selection on data from a customer. Consider a situation in which our method selects the best model for a particular data source. However, we will not stop with this step, but we will rescan for the best model at regular intervals (days, weeks). Using the data thus obtained, we can use unsupervised learning methods to change the model according to the typed scene (day, night, indoor lights on, etc.).

This extension of our method would bring the possibility of automatically changing the object detection model for a specific CCTV/IP camera. We envision that this extension would be mainly applicable when the external environmental conditions change periodically, such as the alternation of natural light and artificial light/night mode of the camera. Additionally, this extension will also bring the ability to adapt to sporadic changes - moving the camera, significant scene change, etc.

We consider improving our object detection method in case of poor results on some videos (e.g., three videos in Table 1). This can also be detected automatically in a deployed situation. The idea is not to process the input data in isolation but as a sequence of frames (perhaps with a higher frame rate, e.g., 0.2 seconds). Thus, we obtain Spatio-temporal data. In this case, we have object identification not only for the current situation but also for the short past and the future (a window of 5 frames would be sufficient). A buffer with object identification from the window frames would serve for additional correction of the current prediction. Some initial experiments were provided, so as a proof of concept, it is promising. We do not present it here. Large areas or inadequate lighting conditions cause some mistakes, so some distant objects are misclassified. Improving lighting conditions, installing more cameras, or requiring object detection only for the closest camera are possibilities for avoiding it. We are working on it and gaining experience.

Semantic analysis of identified objects with respect to the environment could also be an interesting approach for removing false positive predictions. In the case of processing big data from a large number of typed sources (parking lot, office, warehouse, etc.), we could identify the probability of an object's occurrence in a given environment. This approach would be suitable for identifying semantically unacceptable objects in a particular environment - a car in an office, an elephant in a parking lot, etc. Specific knowledge extraction to augment the organizations' databases with information from unstructured CCTV/IP camera data using semantic methods is also left for future work.

## Acknowledgments

This publication was realized with the support of the Operational Program Integrated Infrastructure in the frame of the project: Intelligent technologies for the protection of health-care personnel in the front line and operation of medical facilities during spreading of disease Covid-19, code ITMS2014+: 313011ATQ5 and co-financed by the Europe Regional Development Found.

## References

- [1] Alexey Bochkovskiy, Chien-Yao Wang, Hong-Yuan Mark Liao. (2020) YOLOv4: Optimal Speed and Accuracy of Object Detection, *arXiv.org > cs > arXiv:2004.10934*

- [2] Z. Cai and N. Vasconcelos. (2018) Cascade R-CNN: Delving Into High Quality Object Detection, *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6154–6162, doi: 10.1109/CVPR.2018.00644
- [3] B Cheng, Y Wei, H Shi, R Feris, J Xiong, T Huang. (2020) Decoupled Classification Refinement: Hard False Positive Suppression for Object Detection. *arXiv preprint arXiv:1810.04002*
- [4] Cheng B., Wei Y., Shi H., Feris R., Xiong J., Huang T. (2018) Revisiting RCNN: On Awakening the Classification Power of Faster RCNN. In: Ferrari V., Hebert M., Sminchisescu C., Weiss Y. (eds) In *Computer Vision*, Lecture Notes in Computer Science, vol 11219. Springer, Cham. [https://doi.org/10.1007/978-3-030-01267-0\\_28](https://doi.org/10.1007/978-3-030-01267-0_28)
- [5] M. Feurer, K. Eggensperger, S. Falkner, M. Lindauer, F. Hutter. (2020) Auto-Sklearn 2.0: The Next Generation. *arXiv.org > cs > arXiv:2007.04074*
- [6] Z Ge, Z Song, SX Ding, B Huang. (2017) Data mining and analytics in the process industry: The role of machine learning - *Ieee Access*, - [ieeexplore.ieee.org](https://doi.org/10.1109/ACCESS.2017.2741111)
- [7] R. Girshick, J. Donahue, T. Darrell and J. Malik. (2014) Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 580–587, doi: 10.1109/CVPR.2014.81.
- [8] R. Girshick. (2015) Fast R-CNN, *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1440–1448, doi: 10.1109/ICCV.2015.169.
- [9] Haifeng Jin, Qingquan Song, Xia Hu. (2019) Auto-Keras: An Efficient Neural Architecture Search System. *KDD '19: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Pages 1946–1956 <https://doi.org/10.1145/3292500.3330648>
- [10] T-Y. Lin et al. (2015) Microsoft COCO: Common Objects in Context. *arXiv.org > cs > arXiv:1405.0312*
- [11] Oliff H. , Liu Y. (2017) Towards Industry 4.0 Utilizing Data-Mining Techniques: A Case Study on Quality Improvement, *Procedia CIRP* **63**: 167–172.
- [12] R. Padilla et al. (2020) A Survey on Performance Metrics for Object-Detection Algorithms," *Int. Conf. Systems, Signals and Image Processing (IWSSIP'2020)*, pp. 237–242, doi: 10.1109/IWSSIP48289.2020.9145130, for more see Object-Detection-Metrics on Github
- [13] Padilla, R.; Passos, W.L.; Dias, T.L.B.; Netto, S.L.; da Silva, E.A.B. (2021) A Comparative Analysis of Object Detection Metrics with a Companion Open-Source Toolkit. *Electronics* **10**(3) 279. <https://doi.org/10.3390/electronics10030279>
- [14] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E. (2011), Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research*, **12**, 2825—2830, see also further versions of Scikit-learn on Github
- [15] Vikas C. Raykar, Shipeng Yu, Linda H. Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bogoni, Linda Moy. (2010) Learning From Crowds, *Journal of Machine Learning Research* **11**(43):1297–1322,
- [16] J. Redmon, S. Divvala, R. Girshick and A. Farhadi. (2016) You Only Look Once: Unified, Real-Time Object Detection," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 779–788, doi: 10.1109/CVPR.2016.91.
- [17] S. Ren, K. He, R. Girshick and J. Sun. (2017) Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **39**(6): 1137–1149. doi: 10.1109/TPAMI.2016.2577031.
- [18] Chien-Yao Wang, Alexey Bochkovskiy, Hong-Yuan Mark Liao. (2021) Scaled-YOLOv4: Scaling Cross Stage Partial Network, *arXiv.org > cs > arXiv:2011.08036*
- [19] Z Zou, Z Shi, Y Guo, J Ye. (2019) Object detection in 20 years: A survey, *arXiv preprint arXiv:1905.05055*, 2019 - [arxiv.org](https://arxiv.org/abs/1905.05055)
- [20] Savaglio, Claudio, and Giancarlo Fortino. "A simulation-driven methodology for IoT data mining based on edge computing." *ACM Transactions on Internet Technology (TOIT)* **21.2** (2021): 1–22
- [21] Petrushin, Valery A., and Latifur Khan. *Multimedia data mining and knowledge discovery. Vol. 521. Heidelberg: Springer, 2007*
- [22] Zhu, Xingquan, et al. "Video data mining: Semantic indexing and event detection from the association perspective." *IEEE Transactions on Knowledge and Data engineering* **17.5** (2005): 665–677
- [23] Baughman, Aaron K. and Gao, Jiang and Pan, Jia-Yu and Petrushin, Valery A., *Multimedia Data Mining and Analytics: Disruptive Innovation*, 2016, isbn = 3319347217, Springer Publishing Company, Incorporated
- [24] Tomáš Škopal. On Visualizations in the Role of Universal Data Representation. In *ICMR '20: Proceedings of the 2020 International Conference on Multimedia Retrieval* June 2020 Pages 362–367, <https://doi.org/10.1145/3372278.3390743>
- [25] Tan C., Sun F., Kong T., Zhang W., Yang C., Liu C. (2018) A Survey on Deep Transfer Learning. In: Kůrková V., Manolopoulos Y., Hammer B., Iliadis L., Maglogiannis I. (eds) *Artificial Neural Networks and Machine Learning – ICANN 2018*. ICANN 2018. *Lecture Notes in Computer Science*, vol 11141. pp. 270–279, Springer, Cham. [https://doi.org/10.1007/978-3-030-01424-7\\_27](https://doi.org/10.1007/978-3-030-01424-7_27)
- [26] Petrushin, Valery A. Mining Rare and Frequent Events in Multi-camera Surveillance Video, Pages 80–93, Chapter 5 in [21]
- [27] Wei, Gang (et al.) Multiple-Sensor People Localization in an Office Environment, Pages 427–447, Chapter 21 in [21]
- [28] P. Włodarczyk, J. Soar and M. Ally, "Multimedia data mining using deep learning," *2015 Fifth International Conference on Digital Information Processing and Communications (ICDIPC)*, 2015, pp. 190–196, doi: 10.1109/ICDIPC.2015.7323027.
- [29] A. Nguyen, J. Yosinski and J. Clune, "Deep neural networks are easily fooled: High confidence predictions for unrecognizable images," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 427–436, doi: 10.1109/CVPR.2015.7298640.
- [30] 2020, *Multimedia Big Data Computing for IoT Applications, Concepts, Paradigms and Solutions*, Editors Sudeep Tanwar, Sudhanshu Tyagi, Neeraj Kumar, *Intelligent Systems Reference Library*, vol 163. Springer, Singapore. <https://doi.org/10.1007/978-981-13-8759-3>
- [31] *Multimedia Tools and Applications*, *An International Journal of Springer*